

# Online Neural Network-based Language Identification

Master's Thesis

Daniel H. Draper

INTERACTIVE SYSTEMS LAB - INSTITUTE FOR ANTHROPOMATICS AND ROBOTICS



- 1 Motivation
- 2 Related Work
- 3 Experimental Setup
- 4 DNN Structure
- 5 Post Processing
- 6 Results

- ASR and SLT commercially great success [“OK, Google“, Siri, Jibbigo]
  - Best performing ASR models are still trained on one source language.[1]
- 
- Performance increase for multilingual tasks[1], UI streamlining.

- ASR and SLT commercially great success [“OK, Google“, Siri, Jibbigo]
  - Best performing ASR models are still trained on one source language.[1]
- 
- Performance increase for multilingual tasks[1], UI streamlining.

- ASR and SLT commercially great success [“OK, Google“, Siri, Jibbigo]
  - Best performing ASR models are still trained on one source language.[1]
- 
- Performance increase for multilingual tasks[1], UI streamlining.

## Lecture Translator

- Low latency, online application
- Small source language amount make performance increase likely



## European Parliament

- Attempts to replace translators with SLT [2]
- Many languages challenging, great potential.

## Lecture Translator

- Low latency, online application
- Small source language amount make performance increase likely



## European Parliament

- Attempts to replace translators with SLT [2]
- Many languages challenging, great potential.

## Lecture Translator

- Low latency, online application
- Small source language amount make performance increase likely



## European Parliament

- Attempts to replace translators with SLT [2]
- Many languages challenging, great potential.



## Lecture Translator

- Low latency, online application
- Small source language amount make performance increase likely



## European Parliament

- Attempts to replace translators with SLT [2]
- Many languages challenging, great potential.

## Non-NN-Approaches

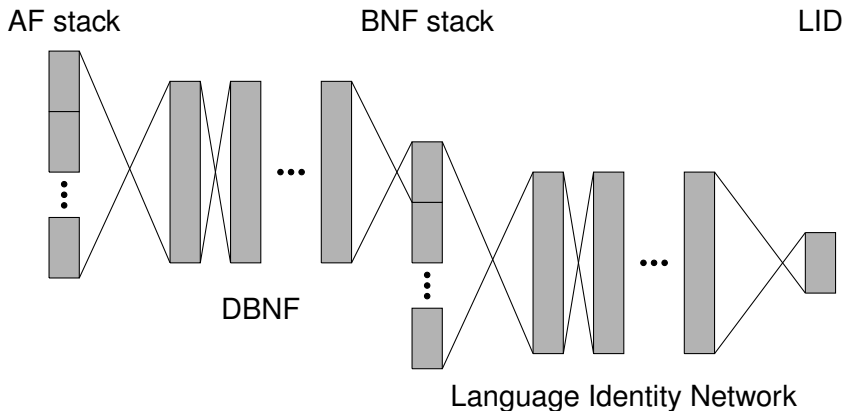
- single-language phone recognition followed by language-dependent, interpolated n-gram language modeling (PRLM)[3]
- multiple single-language phone recognizers and language-dependent parallel phone recognition (PPRM)[3]
- Gaussian Mixture Models (GMM)s [4]
- Vector Space Modelling + PPRM[5]
- I-Vector approaches[6][7]

## Similar Approaches

- Matejka et. al[8] similar setup: BNF  $\rightarrow$  LID, with averaging 5 LID nets.
  - [9] M. Heck et al. evaluate LID approaches prospect of Lecture Translator integration: PPRM/PRLM + Hybrid.
- 
- Markus [10]: "Language Adaptive DNNs for Improved Low Resource Speech Recognition", experimental setup very similar (LFV/LID)

## Similar Approaches

- Matejka et. al[8] similar setup: BNF  $\rightarrow$  LID, with averaging 5 LID nets.
  - [9] M. Heck et al. evaluate LID approaches prospect of Lecture Translator integration: PPRM/PRLM + Hybrid.
- 
- Markus [10]: "Language Adaptive DNNs for Improved Low Resource Speech Recognition", experimental setup very similar (LFV/LID)



## Euronews 2014

- 10 Languages (EN, FR, DE, AR, ES, IT, PO, PT, RU, TR)
- Original Corpus: 72h / Language, Reduced corpus: 18h / Language (Random Sampling of 10.000 speakers)
- 80 % train data, 10 % each dev/test set

## Lecture Data

- 3 Languages (EN, FR, DE)  $\approx$ 10h per language.
- KIT lectures, InterACT25, DGA talks
- Validation of Euronews results, “right” patterns learned (different environment)

## Euronews 2014

- 10 Languages (EN, FR, DE, AR, ES, IT, PO, PT, RU, TR)
- Original Corpus: 72h / Language, Reduced corpus: 18h / Language (Random Sampling of 10.000 speakers)
- 80 % train data, 10 % each dev/test set

## Lecture Data

- 3 Languages (EN, FR, DE)  $\approx$  10h per language.
- KIT lectures, InterACT25, DGA talks
- Validation of Euronews results, “right” patterns learned (different environment)

## European Parliament

- 7 Languages (EN, FR, DE, ES, IT, PO, PT) 3.6h per language
- Simultaneous translation of all parliament speeches
- Further Validation, Proof-of-concept for Lecture Translator integration



## Feature Definition

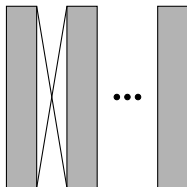
- Samplerate: 16 kHz
- Standard Janus Capabilities:
  - POWER
  - IMEL
  - PITCH
  - TONE
- Context of 6 frames

## Feature Definition

- Samplerate: 16 kHz
- Standard Janus Capabilities:
  - POWER
  - IMEL
  - PITCH
  - TONE
- Context of 6 frames

## Feature Definition

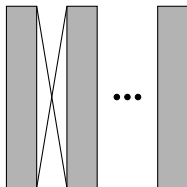
- Samplerate: 16 kHz
- Standard Janus Capabilities:
  - POWER
  - IMEL
  - PITCH
  - TONE
- Context of 6 frames



DBNF

## DBNF

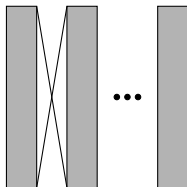
- AF Stack as input
- Targets context-dependent phoneme states
- 6 layers, each 1000 neurons
- BNF of 42 dimensions
- Stacked with context of 11 frames



DBNF

## DBNF

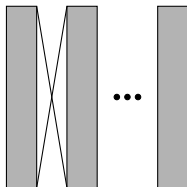
- AF Stack as input
- Targets context-dependent phoneme states
- 6 layers, each 1000 neurons
- BNF of 42 dimensions
- Stacked with context of 11 frames



DBNF

## DBNF

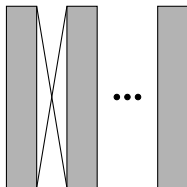
- AF Stack as input
- Targets context-dependent phoneme states
- 6 layers, each 1000 neurons
- BNF of 42 dimensions
- Stacked with context of 11 frames



DBNF

## DBNF

- AF Stack as input
- Targets context-dependent phoneme states
- 6 layers, each 1000 neurons
- BNF of 42 dimensions
- Stacked with context of 11 frames

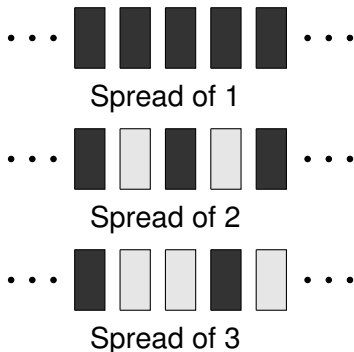


DBNF

## DBNF

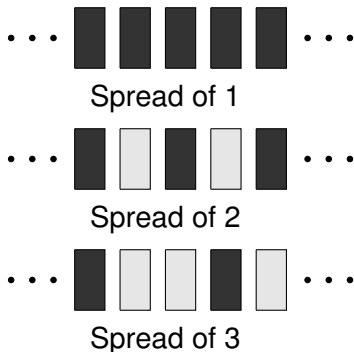
- AF Stack as input
- Targets context-dependent phoneme states
- 6 layers, each 1000 neurons
- BNF of 42 dimensions
- Stacked with context of 11 frames





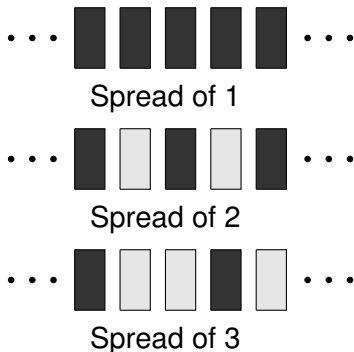
## Spread Evaluation

- In [10], evaluated spreads of 2,3,6 for LFV
- Evaluation of spreads 4 and 5 for different nets on Euronews
- Spread of 3 outperforms 4 and 5, surprisingly 5 outperforms 4



## Spread Evaluation

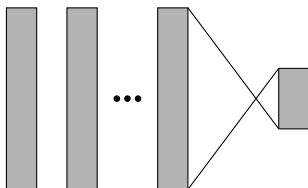
- In [10], evaluated spreads of 2,3,6 for LFV
- Evaluation of spreads 4 and 5 for different nets on Euronews
- Spread of 3 outperforms 4 and 5, surprisingly 5 outperforms 4



## Spread Evaluation

- In [10], evaluated spreads of 2,3,6 for LFV
- Evaluation of spreads 4 and 5 for different nets on Euronews
- Spread of 3 outperforms 4 and 5, surprisingly 5 outperforms 4

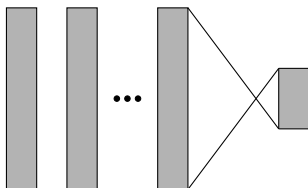
LID



## Baseline Setup

- BNF with context as input
- 5 layers of each 1000 neurons
- Tanh activation, mse loss function
- Mini-batches size 2,000,000, learning rate(lr) 0.01
- Retrained with 1000:10 layer with lr 1

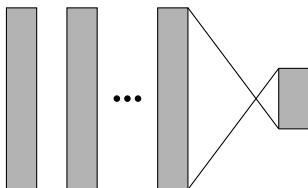
LID



## Baseline Setup

- BNF with context as input
- 5 layers of each 1000 neurons
- Tanh activation, mse loss function
- Mini-batches size 2,000,000, learning rate(lr) 0.01
- Retrained with 1000:10 layer with lr 1

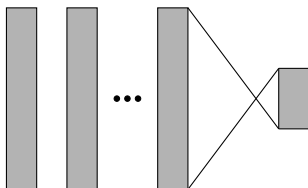
LID



## Baseline Setup

- BNF with context as input
- 5 layers of each 1000 neurons
- Tanh activation, mse loss function
- Mini-batches size 2,000,000, learning rate(lr) 0.01
- Retrained with 1000:10 layer with lr 1

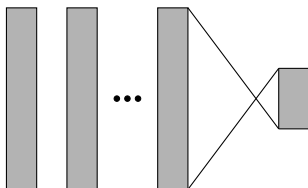
LID



## Baseline Setup

- BNF with context as input
- 5 layers of each 1000 neurons
- Tanh activation, mse loss function
- Mini-batches size 2,000,000, learning rate(lr) 0.01
- Retrained with 1000:10 layer with lr 1

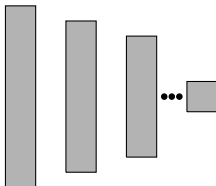
LID



## Baseline Setup

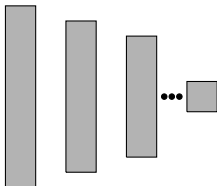
- BNF with context as input
- 5 layers of each 1000 neurons
- Tanh activation, mse loss function
- Mini-batches size 2,000,000, learning rate(lr) 0.01
- Retrained with 1000:10 layer with lr 1





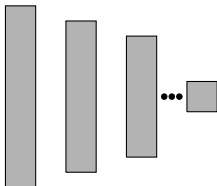
## Other Evaluations

- Lower LR ✓
- Baseline add Layers ✗
- Tree net with 5 layers ✓
- Tree net with 6 layers ✓
- Tree net with > 6 layers ✗ \_
- Cross-Data-Set training ✓
- Concat-Set training ✓✓



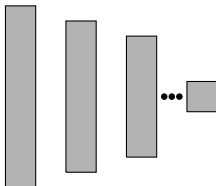
## Other Evaluations

- Lower LR ✓
- Baseline add Layers ✗
- Tree net with 5 layers ✓
- Tree net with 6 layers ✓
- Tree net with > 6 layers ✗ \_
- Cross-Data-Set training ✓
- Concat-Set training ✓✓



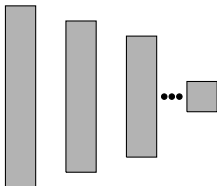
## Other Evaluations

- Lower LR ✓
- Baseline add Layers ✗
- Tree net with 5 layers ✓
- Tree net with 6 layers ✓
- Tree net with > 6 layers ✗ \_
- Cross-Data-Set training ✓
- Concat-Set training ✓✓



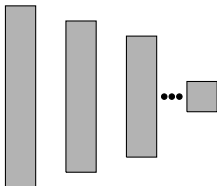
## Other Evaluations

- Lower LR ✓
- Baseline add Layers ✗
- Tree net with 5 layers ✓
- Tree net with 6 layers ✓
- Tree net with > 6 layers ✗ \_
- Cross-Data-Set training ✓
- Concat-Set training ✓✓



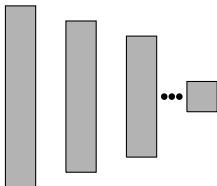
## Other Evaluations

- Lower LR ✓
- Baseline add Layers ✗
- Tree net with 5 layers ✓
- Tree net with 6 layers ✓
- Tree net with  $> 6$  layers ✗
- Cross-Data-Set training ✓
- Concat-Set training ✓✓



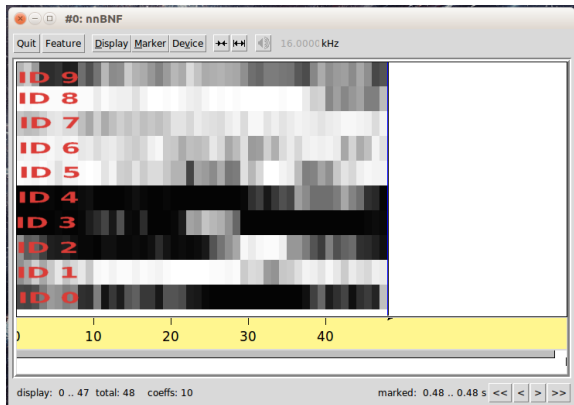
## Other Evaluations

- Lower LR ✓
- Baseline add Layers ✗
- Tree net with 5 layers ✓
- Tree net with 6 layers ✓
- Tree net with  $> 6$  layers ✗ \_
- Cross-Data-Set training ✓
- Concat-Set training ✓✓

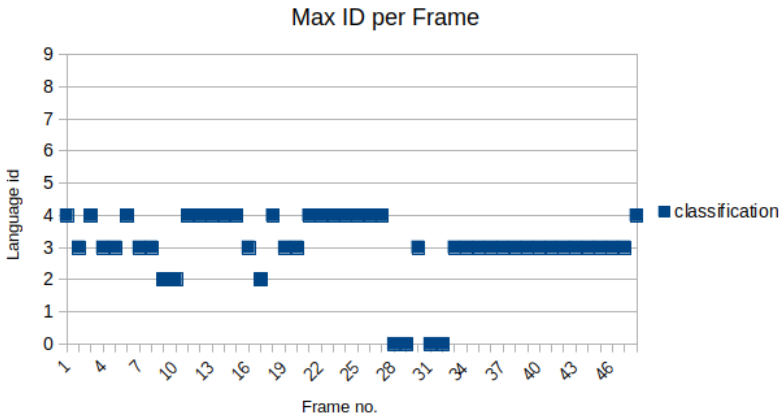


## Other Evaluations

- Lower LR ✓
- Baseline add Layers ✗
- Tree net with 5 layers ✓
- Tree net with 6 layers ✓
- Tree net with  $> 6$  layers ✗ \_
- Cross-Data-Set training ✓
- Concat-Set training ✓✓







## Error Rate

- Count outputs per language per frame
- Max count == actual language → correct, false otherwise
- Metric is percentage of error per language

## Out-Of-Language-Error (OLE)

- Same counting as before
- Sum up FALSELY (as not equal to max language) classified languages
- Metric of “Sureness” of output
- Example: 50 frames classified as DE, 10 as EN, 5 as FR:  
$$(5 + 10) / (5 + 10 + 50) = 0.23$$
- Average OLE per language

## Error Rate

- Count outputs per language per frame
- Max count == actual language → correct, false otherwise
- Metric is percentage of error per language

## Out-Of-Language-Error (OLE)

- Same counting as before
- Sum up FALSELY (as not equal to max language) classified languages
- Metric of “Sureness” of output
- Example: 50 frames classified as DE, 10 as EN, 5 as FR:  
$$(5 + 10) / (5 + 10 + 50) = 0.23$$
- Average OLE per language

## Approaches

- Counting Filter (count in window, in total as maximum only).
- Gaussian Filter (convolution with Gaussian kernel)
- Speech / Noise Filter
- Language Filter (If domain smaller than targets trained)
- Sequence Filter (longest sequence in window)
- Difference (Only count output if difference between 2 max is  $>$  thresh)
- Weighted Average (FILTER capability)

Filter	Overall Error TEST	OLE TEST
Bare Net	0.305	0.198
Gaussian Filter (WS 15)	<b>0.307</b>	0.176
Counting Filter (WS 100)	<b>0.307</b>	<b>0.006</b>
<b>Best</b>	<b>0.307</b>	<b>0.006</b>

- Euronews: Samples  $> 500$  ms: 16 % Error (10L)
- Lecture Data: 6.5 % Error (3L, concat with Euronews data)
- European Parliament: 35 % (7L, 3.5h/language)
- Post Processing: Counting, Gauss Filter promising, relative improvement OLE 97 %.

# Demo (easy)

Play

# Demo (easy)

Language	Number of Frames classified
English	1095
German	123
French	18
<b>Total Error</b>	0.11



# Demo (difficult)

Play

# Demo (difficult)





Language	Number of Frames classified
English	435
German	53
French	99
<b>Total Error</b>	0.25




- RNNs [11]
- Post-Processing promising, further validation or actual focus
- Actual integration into online-environment like LT

- RNNs [11]
- Post-Processing promising, further validation or actual focus
- Actual integration into online-environment like LT



- RNNs [11]
- Post-Processing promising, further validation or actual focus
- Actual integration into online-environment like LT



# Questions?

-  Z. Tang, L. Li, and D. Wang, “Multi-task recurrent model for true multilingual speech recognition,” *arXiv preprint arXiv:1609.08337*, 2016.
-  “Technology and corpora for speech to speech translation,” <http://tcstar.org/>.
-  M. A. Zissman *et al.*, “Comparison of four approaches to automatic language identification of telephone speech,” *IEEE Transactions on speech and audio processing*, vol. 4, no. 1, p. 31, 1996.
-  P. A. Torres-Carrasquillo, E. Singer, M. A. Kohler, R. J. Greene, D. A. Reynolds, and J. R. Deller Jr, “Approaches to language identification using gaussian mixture models and shifted delta cepstral features.” in *Interspeech*, 2002.

-  H. Li, B. Ma, and C. H. Lee, “A vector space modeling approach to spoken language identification,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 15, no. 1, pp. 271–284, Jan 2007.
-  Y. Song, B. Jiang, Y. Bao, S. Wei, and L. R. Dai, “l-vector representation based on bottleneck features for language identification,” *Electronics Letters*, vol. 49, no. 24, pp. 1569–1570, November 2013.
-  L. F. D’haro Enríquez, O. Glembek, O. Plchot, P. Matějka, M. Soufifar, R. d. Córdoba Herralde, and J. Černocký, “Phonotactic language recognition using i-vectors and phoneme posterigram counts,” 2012.



-  P. Matejka, L. Zhang, T. Ng, H. S. Mallidi, O. Glembek, J. Ma, and B. Zhang, “Neural network bottleneck features for language identification,” *Proc. IEEE Odyssey*, pp. 299–304, 2014.
-  M. Heck, S. Ståhlcker, and A. Waibel, “A hybrid phonotactic language identification system with an svm back-end for simultaneous lecture translation,” in *2012 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, March 2012, pp. 4857–4860.

-  M. MÃ¼ller, S. StÃ¼ker, and A. Waibel, “Language adaptive dnns for improved low resource speech recognition,” in *Proceedings of the 17th Annual Conference of the International Speech Communication Association (INTERSPEECH)*, San Francisco, USA, September 8-12 2016.
-  J. Gonzalez-Dominguez, I. Lopez-Moreno, H. Sak, J. Gonzalez-Rodriguez, and P. J. Moreno, “Automatic language identification using long short-term memory recurrent neural networks,” in *Fifteenth Annual Conference of the International Speech Communication Association*, 2014.



H. Li, B. Ma, and K. A. Lee, “Spoken language recognition: From fundamentals to practice,” *Proceedings of the IEEE*, vol. 101, no. 5, pp. 1136–1159, May 2013.

- Given input  $S = \{s(1), s(2), \dots, s(T)\}, \{L_1, L_2, \dots, L_N\}$

## Formulation of LID Task[12]

$$L^O = \arg \max_I p(S|L_I)$$

With segmentation into phones  $v$  (sequence  $\Upsilon$ ):

$$L^O = \arg \max_I P(\Upsilon|L_I)$$

With Viterbi decoding on set of phone models  $M$ :

$$\Upsilon^O = \arg \max_{\Upsilon} P(S|\Upsilon, M)$$

Combining last 2 equations:

$$L^O = \arg \max_I \sum_{\Upsilon} P(S|\Upsilon, M)P(\Upsilon|L_I)$$

- Given input  $S = \{s(1), s(2), \dots, s(T)\}, \{L_1, L_2, \dots, L_N\}$

## Formulation of LID Task[12]

$$L^O = \arg \max_I p(S|L_I)$$

With segmentation into phones  $v$  (sequence  $\Upsilon$ ):

$$L^O = \arg \max_I P(\Upsilon|L_I)$$

With Viterbi decoding on set of phone models  $M$ :

$$\Upsilon^O = \arg \max_{\Upsilon} P(S|\Upsilon, M)$$

Combining last 2 equations:

$$L^O = \arg \max_I \sum_{\Upsilon} P(S|\Upsilon, M) P(\Upsilon|L_I)$$

- Given input  $S = \{s(1), s(2), \dots, s(T)\}, \{L_1, L_2, \dots, L_N\}$

## Formulation of LID Task[12]

$$L^O = \arg \max_I p(S|L_I)$$

With segmentation into phones  $v$  (sequence  $\Upsilon$ ):

$$L^O = \arg \max_I P(\Upsilon|L_I)$$

With Viterbi decoding on set of phone models  $M$ :

$$\Upsilon^O = \arg \max_{\Upsilon} P(S|\Upsilon, M)$$

Combining last 2 equations:

$$L^O = \arg \max_I \sum_{\Upsilon} P(S|\Upsilon, M) P(\Upsilon|L_I)$$

- Given input  $S = \{s(1), s(2), \dots, s(T)\}, \{L_1, L_2, \dots, L_N\}$

## Formulation of LID Task[12]

$$L^O = \arg \max_l p(S|L_l)$$

With segmentation into phones  $v$  (sequence  $\Upsilon$ ):

$$L^O = \arg \max_l P(\Upsilon|L_l)$$

With Viterbi decoding on set of phone models  $M$ :

$$\Upsilon^O = \arg \max_v P(S|v, M)$$

Combining last 2 equations:

$$L^O = \arg \max_l \sum_v \Upsilon P(S|\Upsilon, M)P(\Upsilon|L_l)$$

- Given input  $S = \{s(1), s(2), \dots, s(T)\}$ ,  $\{L_1, L_2, \dots, L_N\}$

## Formulation of LID Task[12]

$$L^O = \arg \max_I p(S|L_I)$$

With segmentation into phones  $v$  (sequence  $\Upsilon$ ):

$$L^O = \arg \max_I P(\Upsilon|L_I)$$

With Viterbi decoding on set of phone models  $M$ :

$$\Upsilon^O = \arg \max_{\Upsilon} P(S|\Upsilon, M)$$

Combining last 2 equations:

$$L^O = \arg \max_I \sum_{\Upsilon} P(S|\Upsilon, M)P(\Upsilon|L_I)$$